

中图法分类号: TP391 文献标识码: A 文章编号: 1006-8961(XXXX)XX-0001-16

论文引用格式: Du Wenliang, Guo Bo, Zhao Jiaqi, Yao Rui, Zhou Yong. XXXX. Semantic-guided contrastive learning for SAR and optical image translation. Journal of Image and Graphics, XX(XX):0001-0016(杜文亮, 郭波, 赵佳琦, 姚睿, 周勇. XXXX. 语义引导对比学习的SAR与光学图像转换. 中国图象图形学报, XX(XX):0001-0016)[DOI:10.11834/jig.250526]

语义引导对比学习的SAR与光学图像转换

杜文亮^{1,2,3}, 郭波^{1,2,3}, 赵佳琦^{1,2,3}, 姚睿^{1,2,3}, 周勇^{1,2,3*}

1. 中国矿业大学计算机科学与技术学院/人工智能学院, 江苏徐州 221116; 2. 矿山数字化教育部工程研究中心, 江苏徐州 221116; 3. 地下空间智能感知与应急物联江苏省产业技术工程化中心, 江苏徐州 221116

摘要: 目的 合成孔径雷达(synthetic aperture radar, SAR)与光学图像转换能够融合两种模态数据的优势, 提供全天候、全天候与高分辨率的观测能力。然而, 当前基于循环一致性生成对抗网络的方法主要侧重于图像结构的宏观重建, 未能充分利用跨模态间的深层语义信息来指导图像生成, 限制了生成图像的语义保真度和在下游任务中的性能。同时, 现有基于对比学习的转换方法在处理遥感图像时, 因同类地物特征高度自相关导致正负样本难以区分, 造成对比机制失效。针对上述问题, 提出了一种语义引导对比学习的SAR与光学图像转换方法。方法 提出了基于语义分割的特征提取模块, 利用预训练的SAR与光学语义分割模型提取像素级语义信息; 提出了语义引导的对比学习模块, 利用先验的语义分割信息, 在对比学习空间中显式构建基于类别一致性的正负样本筛选机制, 有效解决了遥感图像特征同质化导致的传统对比学习失效问题; 设计了融合循环生成结构与对比学习的联合优化框架, 通过引入循环语义分割损失与生成对抗损失, 约束生成图像在结构、纹理和语义层面的一致性。结果 实验在WHU-OPT-SAR和DDHRNet两个公开数据集上进行。实验结果表明, 与当前最优方法相比, 在SAR到光学及光学到SAR的图像转换任务中, 生成质量指标分别最高提升了11.9%和3.8%; 在下游任务中, 语义分割准确率分别提升了16.29%和10.19%, 特征匹配的正确内点比例最高提升了1%。消融实验研究表明, 语义引导对比学习模块与循环语义分割损失对提升模型性能均起到关键作用。结论 本文提出的语义引导对比学习的SAR与光学图像转换方法, 能够有效解决传统对比学习在遥感图像转换中的失效问题, 显著提升了生成图像的语义保真度与跨模态特征对齐能力, 在下游语义分割和图像匹配任务中取得了最优的综合性能, 为无监督SAR与光学图像转换提供了新的解决思路。本文代码开源在链接: <https://www.scidb.cn/s/VVVbnu>。

关键词: 图像转换; 合成孔径雷达; 遥感图像; 对比学习; 语义分割

Semantic-guided contrastive learning for SAR and optical image translation

Du Wenliang^{1,2,3}, Guo Bo^{1,2,3}, Zhao Jiaqi^{1,2,3}, Yao Rui^{1,2,3}, Zhou Yong^{1,2,3*}

1. School of Computer Science and Technology / School of Artificial Intelligence, China University of Mining and Technology, Xuzhou 221116, China; 2. Mine Digitization Engineering Research Center of the Ministry of Education, Xuzhou 221116, China; 3. Jiangsu Provincial Industrial Technology Engineering Center for Intelligent Sensing and Emergency IoT in Underground Space, Xuzhou 221116, China

收稿日期: 2025-10-26; 修回日期: 2026-01-12

* 通信作者: 周勇 yzhou@cumt.edu.cn

基金项目: 国家自然科学基金(No. 62272461, No. 62172417, No. 62002360); 中国矿业大学“双一流”建设提升自主创新能力项目(No. 2022ZZCX06); 澳门科学技术发展基金项目(0020/2024/RIA1)

Supported by: the National Natural Science Foundation of China (Grant No. 62272461, 62172417, and 62002360), and the “Double First-Class” Project of China University of Mining and Technology for Independent Innovation and Social Service (Grant No 2022ZZCX06), and the Science and Technology Development Fund of Macau (Grant No 0020/2024/RIA1).

Abstract: Objective Synthetic Aperture Radar (SAR) and optical remote sensing represent two complementary Earth observation modalities. SAR imagery, thanks to its active imaging mechanism, enables all-weather, all-day observation, while optical imagery provides higher spatial resolution alongside more intuitive visual details. Cross-modal translation between SAR and optical images not only supplements missing data in either modality but also enhances downstream applications such as image matching and semantic segmentation through multi-modal information fusion. However, current methods encounter notable challenges. Approaches based on cycle-consistent generative adversarial networks (CycleGAN) predominantly emphasize macroscopic structural reconstruction and inadequately exploit deep semantic correlations across source and target domains. More critically, traditional contrastive learning techniques face inherent limitations in remote sensing due to the high spatial autocorrelation of similar land-cover features, which blurs the distinction between positive and negative samples, causing the contrastive mechanism to fail. This paper addresses these challenges by proposing a semantic-guided contrastive learning method for SAR and optical image translation, which effectively mitigates feature homogenization in remote sensing image translation and significantly improves semantic fidelity as well as downstream task performance. **Methods** The proposed framework comprises three core components: (1) a semantic feature extraction module, (2) a semantic-guided contrastive learning module, and (3) a joint optimization scheme combining cyclic generation and contrastive learning. Initially, this approach employs pre-trained SAR and optical semantic segmentation models (DeepLabV3) to extract pixel-level semantic features from both real and reconstructed images. Subsequently, the semantic-guided contrastive learning module implements a class-consistency-based positive-negative sample selection strategy: for each query patch in the generated image, patches sharing the same semantic class in the real image serve as positive samples, while those from differing classes are treated as negatives. This strategy effectively mitigates the feature homogenization problem that challenges traditional contrastive learning in remote sensing contexts. Contrastive loss is computed within a shared feature space projected by a lightweight perceptron. Finally, the joint optimization framework integrates cyclic consistency losses at pixel and semantic levels with adversarial losses. The cyclic semantic segmentation loss enforces consistency in structure, texture, and semantics between generated and real images, while adversarial losses enhance image realism. The overall loss function balances these components via weighted hyperparameters. **Results** Comprehensive experiments on two public datasets—WHU-OPT-SAR and DDHRNet—evaluate the proposed method against state-of-the-art approaches, including CycleGAN, contrastive unpaired translation (CUT), query-selected attention (QS-Attn), and conditional diffusion (Con-Diffusion), across image translation quality and downstream tasks. Regarding image translation, our method consistently achieves superior performance. On WHU-OPT-SAR, SAR-to-optical translation yielded a peak signal-to-noise ratio (PSNR) improvement of 11.9% and a mean absolute error (MAE) reduction of 31.1% over the second-best method; optical-to-SAR translation achieved gains of 3.8% in PSNR, 5.6% in structural similarity index measure (SSIM), and a 5.2% reduction in MAE. On DDHRNet, our method sustained leading performance across diverse geographical contexts. For downstream tasks, semantic segmentation and feature matching results confirm marked gains. Semantic segmentation pixel accuracy improved by 16.29% for optical image generation and 10.19% for SAR image generation on WHU-OPT-SAR, outperforming all baselines. For feature matching tasks, the inlier ratio (IR) and euclidean distance-based inlier ratio (IR-ED) improved by up to 1% and 0.49%, respectively. Qualitative comparisons reveal superior visual quality with enhanced detail preservation and grayscale consistency. Ablation studies demonstrate the necessity of each component: removal of the semantic-guided contrastive learning module caused a 10.75% drop in pixel accuracy for optical images, while omitting the cyclic semantic segmentation loss reduced PSNR by 7.0% in optical-to-SAR translation. These findings validate the critical role of our innovations in overcoming remote sensing translation challenges. **Conclusion** The proposed semantic-guided contrastive learning method for SAR and optical image translation, which effectively addresses the prevalent issue of feature homogenization in remote sensing image translation. By integrating semantic segmentation guidance with contrastive learning within a cyclic generative framework, our method substantially enhances the semantic fidelity and downstream utility of generated images, offering a novel solution for unsupervised SAR and optical image translation. Extensive experiments demonstrate that this approach provides a robust and effective solution for unpaired SAR-optical image translation with notable advantages in semantic consistency and cross-modal feature alignment. Our main contributions include: 1) identifying the feature homogenization problem inherent in traditional

contrastive learning for remote sensing translation; 2) proposing a novel semantic-guided contrastive learning framework with category-consistent sample selection; and 3) developing a unified architecture combining cyclic generation and contrastive learning. Future work will explore self-supervised semantic guidance to improve applicability in scenarios lacking semantic annotations, thereby enhancing generalizability and robustness across diverse remote sensing applications. Our codes are available at <https://www.scidb.cn/s/VVVBUu>.

Key words: image-to-image translation; synthetic aperture radar; remote sensing images; contrastive learning; semantic segmentation

0 引言

合成孔径雷达(synthetic aperture radar, SAR)与光学遥感是两种互补的对地观测手段。SAR图像凭借其主动成像机制,能够穿透云雾、不受光照限制,提供全天时、全天候的观测能力。与之相比,光学图像则具有更高的空间分辨率和更符合人类视觉感知的物理特性,信息更直观。因此,研究SAR与光学图像之间的跨模态转换技术,不仅能在任一模态数据缺失时进行补充,更能通过融合多模态信息,为图像匹配、语义分割、图像去雾(Wang等,2025)和目标检测(Hu等,2025)等下游应用提供更鲁棒、更丰富的数据支持。

循环一致性生成对抗网络(cycle-consistent generative adversarial networks, CycleGAN)(Zhu等,2017)首次引入循环一致性损失,解决了非配对数据下的域适应问题。然而,该方法在处理跨模态等风格差异显著的转换任务时,存在模式崩溃的风险。为提升跨模态转换稳定性并支持多属性协同转换,双重对偶生成对抗网络(吴柳玮等,2020)提出通过四个生成器和判别器的对偶架构实现风格与年龄双属性协同转换,并引入专属损失函数保障语义一致性。Torbunov等人(2023)在CycleGAN架构中嵌入视觉Transformer瓶颈层,结合梯度惩罚与自监督预训练缓解模式崩溃问题。在遥感领域,融合残差网络与密集连接卷积网络的并行生成器网络(余佩伦等,2021)提出通过双分支分别处理输入图像,结合基于通道阈值分割的线性插值融合策略优化输出,为跨模态转换的网络结构设计提供了多分支融合的有效思路;Yang等人(2022)以CycleGAN为基础,针对无监督SAR到光学图像转换中生成图像细节不足的问题,提出了包含非平衡生成器、多尺度判别器和综合归一化组的细粒度网络,显著提升了生成图

像的细节质量与物理表征一致性;Guo等人(2024)则针对现有方法中场景记忆不足导致的区域地貌变形和斑点噪声干扰问题,设计了多尺度表示生成器以实现SAR图像场景特征的多尺度融合与利用,构建多感受野判别器以增强不同地貌场景的记忆与生成质量,并引入子带收缩去噪模块有效抑制斑点噪声的影响。然而,当前基于CycleGAN的方法主要侧重于图像结构的宏观重建,未能充分挖掘原始域与目标域之间的深层语义关联。

鉴于此,基于对比无监督转换方法(contrastive unpaired translation, CUT)被提出(Park等,2020),通过引入对比学习机制,充分利用原始图像与生成图像的局部语义信息,通过最大化图像块(Patch)间的互信息实现细节一致性优化,在自然图像转换任务中取得性能提升。基于CUT框架,Kang等人(2020)结合对比损失与条件判别策略,提升转换质量与训练稳定性。为进一步增强对比效果,Wang等人(2021)模型提出实例级困难负样本生成方法,有效提升模型性能;Han等人(2021)通过双编码器结构增强非配对图像间的映射能力。后续研究进一步关注负样本质量与语义结构建模,如查询选择注意力(query-selected attention, QS-Attn)(Hu等,2022)选择关键锚点进行语义对齐,Zhan等人(2022)基于图像对比度动态调整负样本权重,缓解伪影问题;Jun等人(2022)结合语义关系正则化与解耦对比学习,有效提升图像内部异质语义的表达能力和转换多样性。对比学习的核心目标在于最大化正样本对之间的互信息,同时最小化负样本对之间的互信息。然而,遥感图像的特殊性在于:同类地物在特征空间中的分布常高度一致(例如水体和森林等),导致同一语义类别的负样本与正样本在特征空间中的相似度极高甚至趋同,从而使传统对比学习在处理相同语义特征的样本时失效,严重制约了其在遥感图像转换任务中的效能。如图1,蓝、绿和红框分别代表查

询、正样本和负样本图像块。在以 CUT 为代表的传统对比学习方法中,查询图像块与正样本图像块在对比空间中的特征应尽可能相似,而与负样本的特征则应尽可能区分。但在图 1 中,同一语义类别的正样本与负样本的特征极为接近,若继续采用传统对比学习策略,将导致对比机制失效。

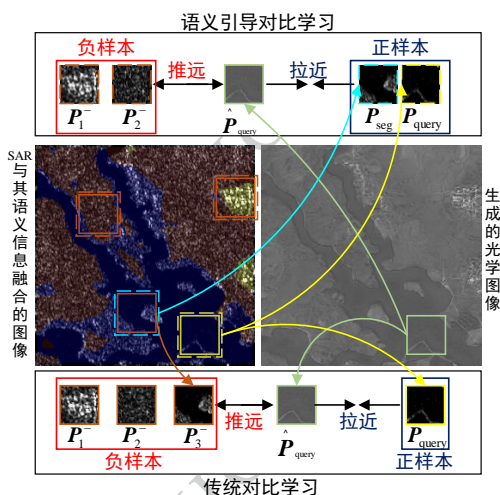


图 1 传统对比学习与语义引导对比学习示意图

Fig. 1 Schematic diagram of traditional contrastive learning and semantic-guided contrastive learning

因此,本文提出语义引导对比学习框架,利用先验的语义分割信息(Jiang 等, 2025),在对比学习空间中显式构建基于类别一致性的正负样本筛选机制(如图 1),以突破遥感图像特征同质化对传统对比学习的约束。同时,为进一步有效利用语义分割信息,在语义引导对比学习框架中加入了 CycleGAN 模型的循环生成结构,并计算循环重构图像与原始图像之间的语义分割损失。

综上,本文主要贡献如下:1)揭示了传统对比学习应用在多源遥感图像转换过程中遇到的特征同质性问题;2)提出语义引导对比学习框架,利用语义分割标签引导查询图像块选择同类区域为正样本,跨类区域为负样本,有效避免遥感图像的同质性问题;3)构建了融合循环生成结构与对比学习的图像转换框架,为非配对遥感图像转换提供了新思路。1 方法

本文提出的语义分割引导对比学习的 SAR 与

光学图像转换方法的整体框架如图 2 所示。该框架主要包括两个核心模块:语义特征提取模块和语义引导的对比学习模块。

如图 2 所示,本方法同时训练了两个生成器: SAR 到光学图像生成器 $G_{S \rightarrow O}$ 和光学到 SAR 图像生成器 $G_{O \rightarrow S}$ 。本文方法将原始的 SAR 图像和光学图像 ($real_{sar}$ 和 $real_{opt}$) 分别输入至生成器 $G_{S \rightarrow O}$ 和 $G_{O \rightarrow S}$, 得到生成的光学和 SAR 图像: $fake_{opt}$ 和 $fake_{sar}$, 然后将 $fake_{opt}$ 和 $fake_{sar}$ 输入至生成器 $G_{O \rightarrow S}$ 和 $G_{S \rightarrow O}$, 得到重建的 SAR 和光学图像: rec_{sar} 和 rec_{opt} 。在生成器训练过程中,原始和重建的图像需输入至语义特征提取模块,获得相应的语义分割图。最后,根据原始和重建的 SAR 与光学图像及其语义分割图,求解分割对比损失、循环语义分割损失、循环损失以及生成对抗损失,以训练生成器和判别器的参数。后续小节,将对语义特征提取模块、语义引导的对比学习模块以及损失函数三个方面进行详细介绍。

1.1 语义特征提取模块

语义特征提取模块由预训练的 SAR 与光学语义分割模型组成。本方法的语义分割模型使用的是 DeepLabV3 (Chen 等, 2017)。其中,在 SAR 图像及其语义分割标注训练集中训练获得的 DeepLabV3 模型用于获得 SAR 图像的语义分割特征,命名为 SAR-Seg; 在光学及其语义分割标注训练集中训练获得的 DeepLabV3 模型用于获得光学图像的语义分割特征,命名为 OPTSeg。

真实的 SAR 与光学图像 ($real_{sar}$ 和 $real_{opt}$) 以及重建的 SAR 与光学图像 (rec_{sar} 和 rec_{opt}) 都需在语义特征提取模块中获得其相应的语义特征。

真实以及重构的 SAR 语义分割特征由相应图像经 SARSeg 分割模型获得,并命名为 Seg_{sar}^{real} 和 Seg_{sar}^{rec} ; 真实以及重构的光学语义分割特征由相应图像经 OPTSeg 分割模型获得,并命名 Seg_{opt}^{real} 和 Seg_{opt}^{rec} 。其中,真实的语义分割特征输入至语义引导的对比学习模块,用于判断正确的正负样本; 真实和重构的 SAR 与光学语义分割特征用于计算生成器的语义损失 \mathcal{L}_{seg} , 约束生成器图像生成的语义一致性, 计算如式(1)所示。式中, \mathcal{L}_{ce} 代表交叉熵损失。

$$\begin{aligned} \mathcal{L}_{seg} &= \mathcal{L}_{ce}(Seg_{sar}^{rec}, \operatorname{argmax}(Seg_{sar}^{real})) + \mathcal{L}_{ce}(Seg_{opt}^{rec}, \operatorname{argmax}(Seg_{opt}^{real})) \end{aligned} \quad (1)$$

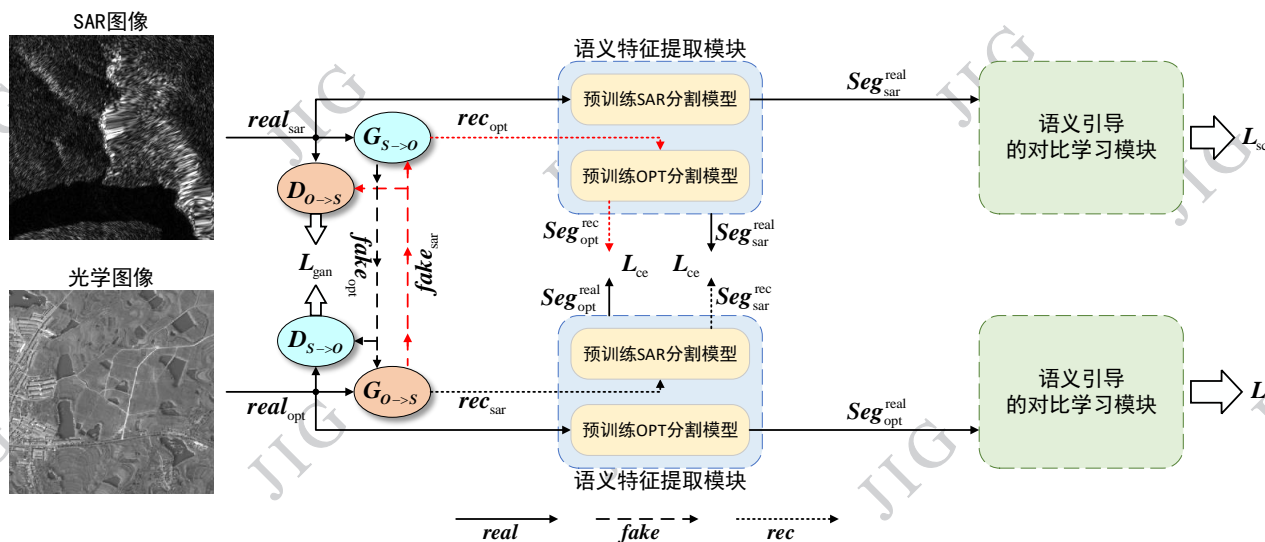


图2 本文提出方法的架构图

Fig. 2 Architecture diagram of the method proposed in this article

1.2 语义引导的对比学习模块

由于 $G_{s \rightarrow o}$ 和 $G_{o \rightarrow s}$ 两个生成器的训练过程都使用了语义引导的对比学习模块,且使用过程一致。因此,本小节以 SAR 图像到光学图像生成为例,介绍语义引导的对比学习模块。语义引导的对比学习模块示意图如图 3 所示,主要包含一个特征编码器 ($E_{s \rightarrow o}^l$) 和一个特征映射感知机 ($MLP_{s \rightarrow o}^l$)。特征编码器 $E_{s \rightarrow o}$ 是 $G_{s \rightarrow o}$ 的一部分,用于提取真实 SAR 图像与生成光学图像的特征。 $E_{s \rightarrow o}^l$ 表示编码器 $E_{s \rightarrow o}$ 第 l 层的输出。其中, $l \in \{1, 2, 3, \dots, L\}$, L 表示 $E_{s \rightarrow o}$ 为从 $G_{s \rightarrow o}$ 提取的 L 层网络。在网络架构层面,本模块沿用了 CUT 模型 (Park 等, 2020) 中提出的图像块级对比学习框架,通过特征编码器 ($E_{s \rightarrow o}$) 提取 L 层特征,并利用与 CUT 相同结构的多层感知机 (multilayer perceptron, MLP) 将特征映射至对比空间。然而,传统的 CUT 方法在对比学习策略上仅基于空间位置划分正负样本,即同一位置的图像块为正样本其余为负样本。这种策略在处理具有高纹理重复性的遥感图像时容易失效。与之不同,本文提出的语义引导对比学习模块在对比学习架构的基础上,设计并引入了语义一致性筛选机制,以确保正负样本的语义一致性,其具体过程如下。

在构建语义引导的对比学习正样本与负样本过程中,首先在生成图像 ($fake_{opt}$ 或 $fake_{sar}$) 中选取的一个图像块 (Patch) 作为生成查询块 \hat{p}_{query} ,并将真实图像中对应位置的图像块作为真实查询块 p_{query} 。然

后,根据 p_{query} 在真实图像语义分割特征 Seg_{sar}^{real} 中对应的类别,获取真实图像中同类别的图像块 p_{seg} ,并将不同类别的图像块作为负样本图像块 p_- 。正样本图像块 p_+ 则为同类别图像块 p_{seg} 与真实查询块 p_{query} 的并集,即 $p_+ = \{p_{query}, p_{seg}\}$ 。传统对比学习方法通常采用 1×1 像素作为图像块的尺寸,此时 p_+ 中图像块的个数与 p_- 中图像块个数之和为图像的像素个数。

接下来,利用特征编码器 $E_{s \rightarrow o}^l$ 和特征映射感知机 $MLP_{s \rightarrow o}^l$ 将 \hat{p}_{query} 、 p_{query} 、 p_{seg} 、 p_+ 和 p_- 映射到对比空间中,得到相应的 C_l 维的特征向量: $\hat{v}_{query}^l \in \mathbf{R}^{C_l}$ 、 $v_{query}^l \in \mathbf{R}^{C_l}$ 、 $v_{seg}^l \in \mathbf{R}^{P \times C_l}$ 、 $v_+^l \in \mathbf{R}^{(P+1) \times C_l}$ 和 $v_-^l \in \mathbf{R}^{N \times C_l}$ 。其中, P 为相同类别块的个数, N 为不同类别图像块的个数, C_l 为特征编码器第 l 层特征维度。设置语义引导的交叉熵损失,令 \hat{v}_{query}^l 与 v_+ 更相似、与 v_- 更不同。语义引导的交叉熵损失如式 (2) 所示:

$$\begin{aligned} & \ell(\hat{v}_{query}^l, v_+^l, v_-^l) \\ &= -\log \left[\frac{\sum_{v_+^l \in Posi} \exp(\hat{v}_{query}^l \cdot v_+^l / \tau)}{\sum_{v_+^l \in Posi} \exp(\hat{v}_{query}^l \cdot v_+^l / \tau) + \sum_{v_-^l \in Neg} \exp(\hat{v}_{query}^l \cdot v_-^l / \tau)} \right] \end{aligned} \quad (2)$$

式中, $Posi$ 和 Neg 分别代表正样本和负样本特征向量的集合, $Posi = \{v_{query}^l, v_{seg}^{l(1)}, \dots, v_{seg}^{l(S)}\}$, $Neg = \{v_-^{l(1)}, \dots, v_-^{l(N)}\}$, τ 表示缩放系数。

由于 SAR 到光学和光学到 SAR 生成器的训练过程都使用了语义引导的对比学习模块训练,因此编码器 $E_{s \rightarrow o}$ 和 $E_{o \rightarrow s}$ 在选取的 L 个特征层上都计算了

两个模态在对比空间中的对比损失,以确保两个模态特征在对比空间的语义一致性。同时,还需在两编码器对应的 L 个特征层计算各模态的一致性损失(Identity Loss),以约束各模态在对比空间中语义的不变性。其中,对于编码器的第 l 层,查询块共有 Q_l 个。综上,语义引导的对比学习损失 \mathcal{L}_{sc} 的计算过程如式(3)所示:

$$\begin{aligned} \mathcal{L}_{sc} &= \mathbb{E}_{s \sim S} \sum_{l=1}^L \sum_{q=1}^{Q_l} \ell(\hat{\mathbf{v}}_{\text{query}}^{l(q)}, \mathbf{v}_+^{l(q)}, \mathbf{v}_-^{l(q)}) \\ &+ \mathbb{E}_{o \sim O} \sum_{l=1}^L \sum_{q=1}^{Q_l} \ell(\hat{\mathbf{v}}_{\text{query}}^{l(q)}, \mathbf{v}_+^{l(q)}, \mathbf{v}_-^{l(q)}) \\ &+ \mathbb{E}_{s \sim S} \sum_{l=1}^L \sum_{q=1}^{Q_l} \ell(\hat{\mathbf{v}}_{\text{query}}^{l(q)}, \mathbf{v}_+^{l(q)}, \mathbf{v}_-^{l(q)}) \\ &+ \mathbb{E}_{o \sim O} \sum_{l=1}^L \sum_{q=1}^{Q_l} \ell(\hat{\mathbf{v}}_{\text{query}}^{l(q)}, \mathbf{v}_+^{l(q)}, \mathbf{v}_-^{l(q)}) \end{aligned} \quad (3)$$

式中, \mathbb{E} 表示数学期望, S 和 O 分别代表真实的SAR和光学图像的数据集, s 和 o 分别代表数据集中真实的SAR和光学图像, $\hat{\mathbf{v}}_{\text{query}}^l$ 表示一致性生成图像(idt_{sar} 或 idt_{opt})查询块对应的特征向量, idt_{sar} 和 idt_{opt} 由 $G_{O \rightarrow S}$ (real_{sar})和 $G_{S \rightarrow O}$ (real_{opt})得到。式(3)的前半部分为两模态间在对比空间的对比损失,后半部分为各模态内在对比空间的一致性损失。

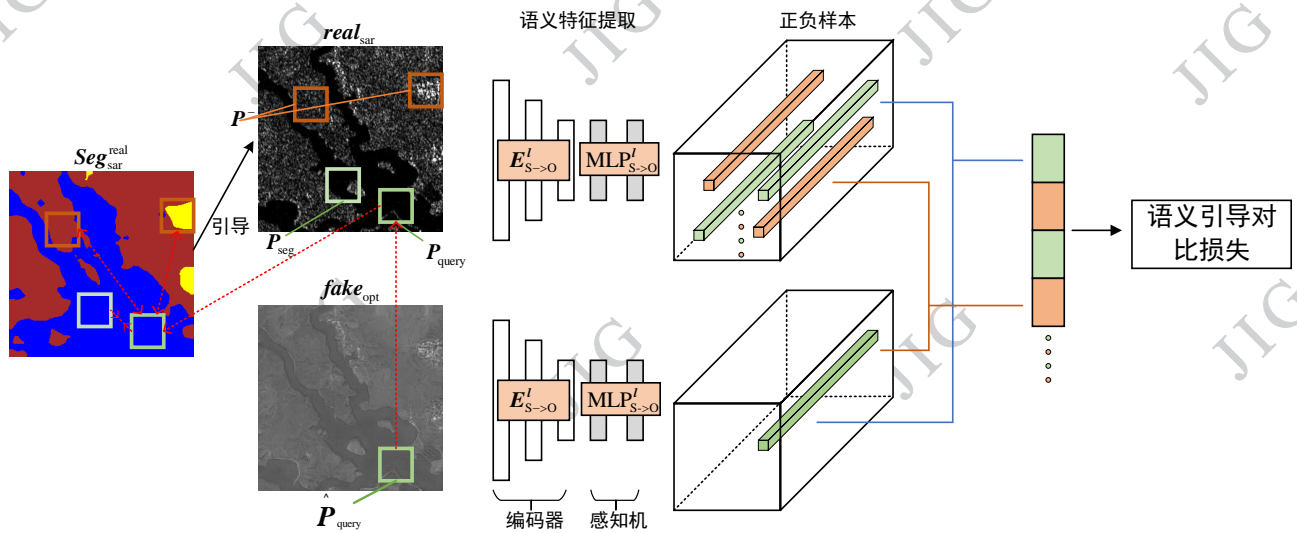


图3 语义引导的对比学习模块。 \hat{p}_{query} 为生成查询向量; p_{query} 为真实查询向量; p_{seg} 为与查询向量相同类别的特征向量; p 为负样本特征向量

Fig. 3 Semantic-guided contrastive learning module. \hat{p}_{query} is the generated query vector; p_{query} is the true query vector; p_{seg} is the feature vector of the same category as the query vector; p is the negative sample feature vector

1.3 损失函数

本文方法除了使用了前两小节提到的语义损失 \mathcal{L}_{seg} 和语义引导的对比学习损失 \mathcal{L}_{sc} ,还使用了生成

对抗损失和循环损失以保证SAR到光学基础的图像生成质量。

生成对抗损失如式(4)所示:

$$\begin{aligned} \mathcal{L}_{\text{gan}} &= \mathbb{E}_{o \sim O} \log D_{S \rightarrow O}(\text{real}_{\text{opt}}) + \mathbb{E}_{s \sim S} \log(1 - D_{S \rightarrow O}(G_{S \rightarrow O}(\text{real}_{\text{sar}}))) \\ &+ \mathbb{E}_{s \sim S} \log D_{O \rightarrow S}(\text{real}_{\text{sar}}) + \mathbb{E}_{o \sim O} \log(1 - D_{O \rightarrow S}(G_{O \rightarrow S}(\text{real}_{\text{opt}}))) \end{aligned} \quad (4)$$

循环损失函数计算过程如式(5)所示。式中, λ_1 和 λ_2 为超参数,分别设置为10和0.5。

$$\begin{aligned} \mathcal{L}_{\text{cyc}} &= \lambda_1 \left(\left\| \text{real}_{\text{sar}} - \text{rec}_{\text{sar}} \right\|_1 + \left\| \text{real}_{\text{opt}} - \text{rec}_{\text{opt}} \right\|_1 \right) \\ &+ \lambda_2 \left(\left\| \text{real}_{\text{sar}} - \text{idt}_{\text{sar}} \right\|_1 + \left\| \text{real}_{\text{opt}} - \text{idt}_{\text{opt}} \right\|_1 \right) \end{aligned} \quad (5)$$

本文总损失函数为 \mathcal{L}_{all} 。式中, λ_3 、 λ_4 和 λ_5 为超参数, 统一设置为 1。

$$\mathcal{L}_{\text{all}} = \mathcal{L}_{\text{cyc}} + \lambda_3 \mathcal{L}_{\text{gan}} + \lambda_4 \mathcal{L}_{\text{seg}} + \lambda_5 \mathcal{L}_{\text{sc}} \quad (6)$$

2 实验结果及分析

2.1 对比方法与实验设置

本文对比方法包括两个基线方法 CycleGAN 和 CUT 以及当前最优的基于对比学习的图像转换方法 QS-Attn 和基于扩散模型的图像转换方法条件扩散 (conditional diffusion, Con-Diffusion) (Bai 等, 2023)。

本文实验均在包含 1 个 NVIDIA GeForce RTX 4090D GPU 的服务器上完成。批次大小设置为 2, 训练 epoch 数设置为 100, 优化器采用 Adam。学习率在前 50 个 epoch 训练过程中设置为 2×10^{-4} 在后 50 个 epoch 以线性衰减到 0。式(2)中 τ 设置为 0.07。

2.2 数据集

实验采用两个 SAR 与光学图像语义分割数据集: WHU-OPT-SAR (Li 等, 2022) 和 DDHRNet (Ren 等, 2022)。WHU-OPT-SAR 数据集由高分一号拍摄的光学图像与高分三号拍摄的 SAR 图像构成。该数据集的训练集、验证集和测试集分别包含 20580 对、5880 对和 2940 对分辨率为 256×256 的 SAR 和光学图像及其对应的语义分割标注。由于现有图像转换模型训练速度较慢, 本文从训练集中随机选取 2000 对图像, 用于训练本文提出的模型及四种对比方法。

DDHRNet 数据集由高分二号拍摄的光学图像与高分三号拍摄的 SAR 图像构成, 覆盖中国西安、山东及韩国区域, 本文实验仅使用山东区域数据。该数据集经随机采样后按 8:1:1 的比例划分为训练集、验证集和测试集, 分别包含 6264 对、783 对和 784 对分辨率为 256×256 的 SAR 与光学图像及其语义分割标注。

两个数据集中的 SAR 与光学图像均为严格配对 (Paired), 可用于有效评估图像转换质量。然而, 为模拟仅有非配对数据时的训练场景, 本文在训练图像转换模型时, 对 SAR 与光学图像分别进行随机

采样, 使得每个批次中的图像均为非配对 (Unpaired) 形式。

2.3 评估指标

本文从图像转换质量与下游任务性能两个维度对比模型的 SAR 到光学图像转换质量进行评价。其中, 图像转换质量评估采用峰值信噪比 (peak signal-to-noise ratio, PSNR)、结构相似性指标 (structural similarity index measure, SSIM) 以及平均绝对误差 (mean absolute error, MAE) 作为评估指标。

下游任务评估包括语义分割与图像匹配两项任务。在语义分割任务中, 本文采用在语义特征提取模块中针对 SAR 和光学图像分别训练的 DeepLabV3 模型 (SARSeg 和 OPTSeg), 对转换后的图像进行语义分割, 并使用像素准确率 (pixel accuracy, PA) (Badrinarayanan 等, 2017) 对分割结果进行评估。

在特征匹配任务中, 使用在两数据集上训练的特征匹配模型 (Yagmur 等, 2024) 进行特征点提取与粗匹配, 并将粗匹配点数量与提取特征点数量的比值 (putative matches ratio, PMR) 作为第一个评价指标。然后使用随机抽样一致 (random sample consensus, RANSAC) 算法筛选出正确匹配点, 以正确匹配点数量占粗匹配点数量的比例 (inlier ratio, IR) 作为第二个评估指标。此外, 由于两个数据集中的 SAR 与光学图像均为严格配对, 可通过计算匹配特征点的欧氏距离直接判定正确匹配点, 将欧氏距离小于 3 个像素的匹配点定义为欧氏距离正确匹配, 并将其数量与粗匹配点数量的比值 (inlier ratio of euclidean distance, IR-ED) 作为第三个评估指标。

2.4 对比实验和分析

2.4.1 定量实验结果与分析

本小节首先对各类方法在图像转换质量方面的定量实验结果进行对比分析, 进而评估其在下游任务中的表现。

各方法在 WHU-OPT-SAR 和 DDHRNet 数据集上的 SAR 与光学图像跨模态转换任务的图像质量定量实验结果如表 1 和表 2 所示。

在两个数据集上, 本文方法与 CycleGAN 在 SAR 到光学图像生成及光学到 SAR 图像生成的双向任务中, 均展现出优于其他对比方法的图像生成质量。因此, 若仅从传统的图像生成质量评估指标来看, 基于循环生成方法在 SAR 与光学图像的跨模态生成中, 与真实图像具有更高的像素相似度 (如 PSNR、

MAE)和结构相似性(如SSIM)。相比之下,基于对比学习的方法CUT和QS-Attn在图像生成质量评估指标上的表现弱于CycleGAN和本文方法。此外,采用查询选择注意力模块的QS-Attn方法相比基线方法CUT,在SAR与光学图像转换任务中并未表现出显著优势。而基于扩散模型的Con-Diffusion方法在两个数据集上的图像生成质量评估指标显著弱于其他方法,表明扩散模型在SAR与光学图像转换任务的图像质量评估方面不具备优势。

具体而言,在WHU-OPT-SAR数据集上的SAR到光学图像生成任务中,本文方法在PSNR和MAE两项指标上均达到最优性能,相较于次优方法CycleGAN和CUT,PSNR提升了2.2464(11.9%),MAE(越低越好)降低了0.0345(31.1%);而在SSIM指标上,本方法较CycleGAN降低0.045(10.6%)。在光学图像到SAR图像生成任务中,本文方法在所有三项图像质量评价指标上均取得最优表现,相较次优方法在PSNR和SSIM分别提升了0.4771(3.8%)和0.003(5.6%),在MAE降低了0.0094(5.2%)。

在DDHRNet数据集上,本文方法在SAR到光学图像生成任务中的三项图像质量评价指标(PSNR、SSIM、MAE)上均取得最优表现,相较于次优方法CycleGAN,PSNR提升了0.2524(1.6%),SSIM提升了0.0216(8.9%),MAE降低了0.0071(5.5%)。然而,在光学到SAR图像生成任务中,CycleGAN在三项指标上均表现最优,而本文方法未能达到次优水平。这表明,在CycleGAN基础上引入语义引导的对比学习机制,在DDHRNet数据集上生成的SAR图像在图像质量方面并未产生预期提升效果。

表3展示了各类方法在WHU-OPT-SAR数据集上,图像生成结果在语义分割(PA)和特征匹配(PMR、IR和IE-ED)两个下游任务的实验结果。本文提出方法在两数据集的两个图像转换任务上都取得了最优结果。具体而言,本文方法在光学图像生成与SAR图像生成的语义分割结果分别比次优方法分别高出16.29%和10.19%。CycleGAN在WHU-OPT-SAR数据集上的光学生成任务中未取得次优结果,但在其余三个实验中均表现次优。这表明,基于循环生成方法在语义分割下游任务中仍占据显著优势。值得一提的是,基于扩散模型的Con-Diffusion方法在语义分割任务中并未像其在图像质

量评估中那样表现出明显劣势,甚至在SAR生成图像的语义分割结果中,优于两种基于对比学习的方法。

在特征匹配下游任务中,此前在图像生成质量评估中表现最弱的Con-Diffusion方法,在PMR指标上却取得了最优结果。这表明,尽管其生成图像在像素级相似度上与真实图像存在较大差距,但在特征级相似度上表现较强,从而使其生成的59.87%特征点能够构成粗匹配。基于对比学习的方法同样取得了优异的特征匹配结果:CUT方法虽在PMR指标上比Con-Diffusion方法略低1.78%,但其精匹配比例比Con-Diffusion方法高2.73%(IR)和0.49%(IR-ED)。与之相对,在图像生成质量评估中表现优异的CycleGAN,在特征匹配实验中表现最差,仅有12%的特征点可作为粗匹配,且在这些粗匹配中,仅有2%(IR)和0.2%(IR-ED)的特征点能够通过RANSAC方法和欧氏距离筛选成为精匹配。本文方法在PMR指标上取得了次优结果,并在精匹配中均取得了最优性能,表明本文方法有效结合了CycleGAN的图像生成质量与CUT的特征生成质量优势。

2.4.2 定性实验结果与分析

本小节首先对各类方法的图像转换案例进行定性分析,随后对下游语义分割任务和特征匹配任务的案例进行定性分析。所有定性分析均基于WHU-OPT-SAR数据集的案例展开。

图4(a)和(b)分别展示了在WHU-OPT-SAR数据集上,各方法在SAR到光学图像和光学图像到SAR转换任务中的4组生成结果,其中左起第一列分别为真实的SAR和光学图像。图像转换结果进一步印证了前一小节定量实验的结论:从像素级图像转换结果来看,基于循环生成的方法(CycleGAN和本文方法)的转换效果优于基于对比学习的方法(CUT和QS-Attn)以及基于扩散模型的方法(Con-Diffusion)。

具体而言,Con-Diffusion生成的图像在像素级上与真实图像存在显著差异,尤其是在光学图像转换SAR图像方面,其生成SAR图像的灰度值与真实SAR图像存在明显偏差。例如,在图4(b)的第一和第二行光学到SAR图像生成结果中,Con-Diffusion生成的SAR图像整体呈现灰白色调,而真实SAR图像则以灰黑色调为主。CUT和QS-Attn方法生成的

表 1 各方法在 WHU-OPT-SAR 数据集上进行图像转换的定量实验结果

Table 1 Quantitative experimental results of image translation using various methods on the WHU-OPT-SAR dataset

方法	SAR→OPT			OPT→SAR		
	PSNR↑	SSIM↑	MAE↓	PSNR↑	SSIM↑	MAE↓
CycleGAN(Zhu 等, 2017)	<u>18.8835</u>	0.4243	0.1135	<u>12.4394</u>	<u>0.0532</u>	<u>0.1799</u>
CUT(Park 等, 2020)	18.5940	0.3599	<u>0.1110</u>	11.7418	0.0416	0.1933
QS-Attn(Hu 等, 2022)	18.4340	0.3682	0.1125	11.7803	0.0414	0.1925
Con-Diffusion(Bai 等, 2023)	14.9459	0.2455	0.1922	8.1529	0.0208	0.3306
本文方法	21.1299	<u>0.3793</u>	0.0765	12.9165	0.0562	0.1705

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

表 2 各方法在 DDHRNet 数据集上进行图像转换的定量实验结果

Table 2 Quantitative experimental results of image translation using various methods on the DDHRNet dataset

方法	SAR→OPT			OPT→SAR		
	PSNR↑	SSIM↑	MAE↓	PSNR↑	SSIM↑	MAE↓
CycleGAN(Zhu 等, 2017)	<u>16.2501</u>	<u>0.2423</u>	<u>0.1285</u>	11.8008	0.0998	0.1981
CUT(Park 等, 2020)	15.3855	0.2222	0.1358	11.5565	<u>0.0914</u>	0.2039
QS-Attn(Hu 等, 2022)	15.4356	0.2251	0.1346	<u>11.6227</u>	0.0906	<u>0.2021</u>
Con-Diffusion(Bai 等, 2023)	10.1888	0.1218	0.2935	9.8624	0.0814	0.2595
本文方法	16.5025	0.2639	0.1214	11.5337	0.0837	0.2025

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

表 3 各方法在 WHU-OPT-SAR 数据集的语义分割和特征匹配实验结果

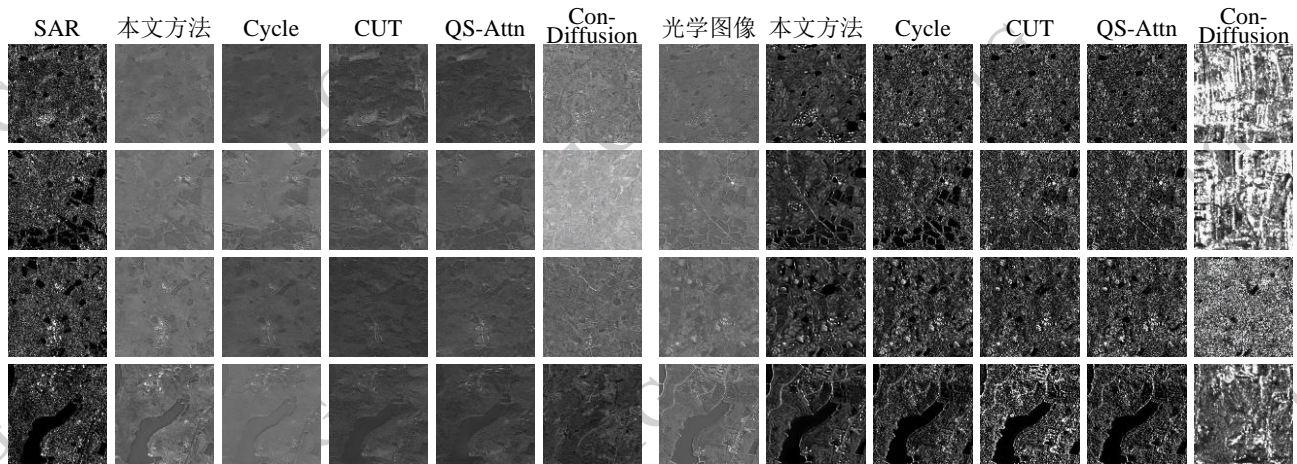
Table 3 Experimental results of semantic segmentation and feature matching for various methods on the WHU-OPT-SAR dataset

方法	PA(SAR→OPT)↑	PA(OPT→SAR)↑	PMR↑	IR↑	IR-ED↑
CycleGAN(Zhu 等, 2017)	<u>0.4709</u>	<u>0.4042</u>	0.1200	0.0206	0.0020
CUT(Park 等, 2020)	0.3728	0.4367	0.5809	<u>0.3172</u>	<u>0.1350</u>
QS-Attn(Hu 等, 2022)	0.3601	0.4410	0.5086	0.3137	0.0775
Con-Diffusion(Bai 等, 2023)	0.4202	<u>0.4576</u>	0.5987	0.2899	0.1301
本文方法	0.4961	0.5061	<u>0.5886</u>	0.3271	0.1390

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好。

图像整体质量优于 Con-Diffusion,但在 SAR 到光学图像生成任务中,部分区域仍存在模糊现象。例如,在图 4(a)的第四行 SAR 到光学图像生成案例中,CUT 和 QS-Attn 方法生成的光学图像上半部分的地

表纹理与建筑边缘细节生成结果较为模糊。相比之下,本文方法与 CycleGAN 在该区域生成的图像细节更为丰富。此外,本文方法生成图像在与真实图像的灰度一致性方面,显著优于其他对比方法。例如,



(a) SAR转换为光学图像的实验效果 (b) 光学图像转换为SAR图像的实验效果

((a) SAR-to-optical image translation results; (b) Optical-to-SAR image translation results)

图4 各方法在WHU-OPT-SAR数据集上的转换效果

Fig. 4 Qualitative comparison of different methods on the WHU-OPT-SAR dataset

在图4的第一和第三行SAR到光学图像生成案例中,本文方法生成的光学图像与真实图像在灰度上具有高度一致性,而其他方法生成的图像整体灰度值较真实图像偏暗。

图5(a)和(b)分别展示了各方法在图4(a)和(b)上生成的4组图像对应的语义分割结果,其中,左起三列依次为真实图像、语义分割标注,以及分割器OPTSeg和SARSeg对真实图像的分割结果。与定量实验结论一致,在语义分割引导的加持下,本文提出的方法在语义分割下游任务中的表现显著优于其他方法。

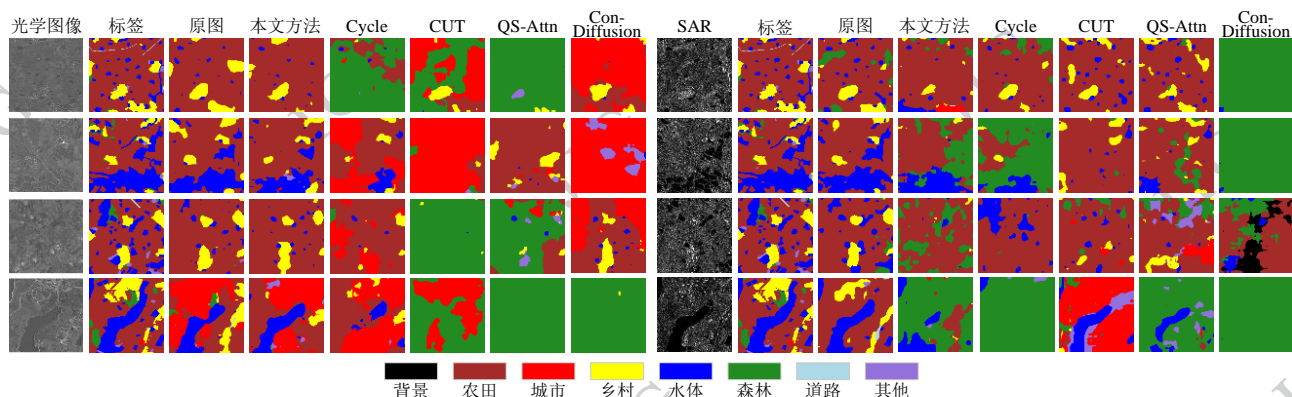
具体而言,在图5(a)所示的生成光学图像的语义分割结果中,基于对比学习的方法中仅有QS-Attn方法在第二行中生成的农田特征,以及CUT方法在第一行中生成的部分乡村特征与真实图像较为一致,使其语义分割结果与标注较为接近。其余生成光学图像在OPTSeg分割器下的表现与真实图像分割结果以及分割标注的差异较为显著。例如,CUT方法在四组图像中,将原本为农田的特征生成成为城市或森林;QS-Attn方法则将第一、第三和第四行中的农田及水体特征生成成为森林,导致语义分割出现大面积错误。相比之下,基于扩散模型的Con-Diffusion和基于循环生成的CycleGAN的生成效果略优,尤其在第三行中,两种方法生成的农田特征与真实图像较为一致。然而,CycleGAN和Con-Diffusion在第一行中分别将农田特征生成成为森林和城市,且Con-Diffusion在第四行中生成的光学图像

特征更接近森林,导致整幅图像分割错误。与此相对,本文方法在四组光学图像生成中的特征与真实图像更为接近。尽管其在第四行中的分割结果将农田和乡村的组合误分割为城市,这主要是由于第四行中的特征更为复杂,导致OPTSeg分割器将大片农田与乡村的组合识别为城市。

图6展示了SAR图像转换至光学图像后与真实图像的特征匹配可视化结果。该结果与前述定量分析的结论相吻合。具体而言,CycleGAN虽然在像素级的图像转换上取得了较好的视觉效果,但在特征匹配任务中表现不佳,在全部四个测试案例中,仅有一对实现了正确匹配。与之形成对比的是,Con-Diffusion方法,尽管其图像转换质量相对较低,却在特征匹配上表现出色,其生成的粗匹配对数量(图中绿色与红色线的总和)在所有对比方法中达到最多。基于对比学习的CUT和QS-Attn方法,虽然粗匹配总数的绝对值不高,但其中正确匹配(图中绿色直线)的占比较高。本文所提方法则兼具二者优势,不仅获得了数量可观的粗匹配对,同时保持了较高的正确匹配比例。综上,本文方法在像素级图像生成与特征级匹配两个任务上,均展现出卓越的综合性能。

2.5 消融实验

本文提出方法相较于仅采用生成对抗损失 \mathcal{L}_{gan} 的基线模型,主要引入了循环一致性架构及其对应损失 \mathcal{L}_{cyc} 、语义特征提取模块及其语义分割损失 \mathcal{L}_{seg} 和语义引导的对比学习模块及其对应损失 \mathcal{L}_{sc} 。为

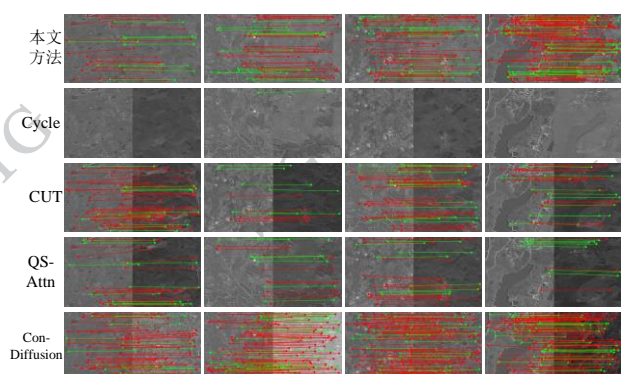


(a) SAR 转换到光学图像后的分割效果 (b) 光学图像转换到 SAR 后的分割效果

((a) Segmentation results after SAR-to-optical image translation; (b) Segmentation results after optical-to-SAR image translation)

图 5 各方法在 WHU-OPT-SAR 数据集转换后的分割效果

Fig. 5 Segmentation results of different methods after image translation on the WHU-OPT-SAR dataset



绿色和红色线段分别表示 RANSAC 识别的正确匹配和错误匹配

图 6 各方法在 WHU-OPT-SAR 数据集的配准效果

Fig. 6 Registration performance of various methods on the WHU-OPT-SAR dataset

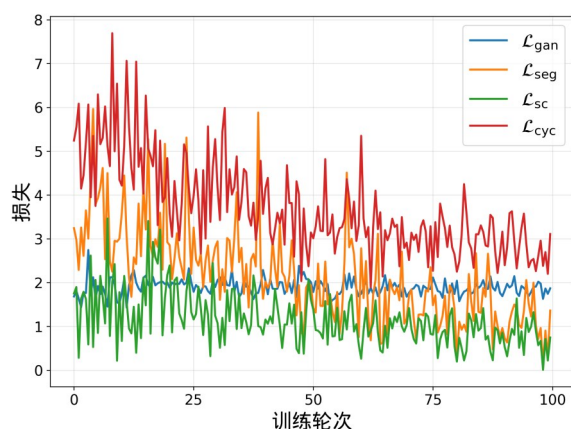


图 7 所提方法各项损失在训练过程中的变化曲线

Fig. 7 Training loss curves of the proposed method

验证各模块的有效性及其损失函数对模型性能的独立贡献,在 WHU-OPT-SAR 数据集上设计了逐步引

入上述损失函数的消融实验。该实验从图像转换质量及下游语义分割准确性两个维度评估模型表现。表 4 展示了不同损失函数组合下的实验结果。

实验结果表明,在仅使用生成对抗损失 \mathcal{L}_{gan} 时,由于缺乏有效的约束,SAR 与光学图像的转换质量较差,各项指标均处于最低水平。在此基础上引入循环一致性损失 \mathcal{L}_{cyc} 后(即等同于 CycleGAN 基线模型),模型的结构保持能力显著增强,SAR 到光学转换的 PSNR 提升至 18.2425,初步实现了跨模态转换。

为了验证语义信息的引入对转换性能的提升,在基线模型上逐步叠加了不同的语义相关损失项。当增加语义分割损失($\mathcal{L}_{gan} + \mathcal{L}_{cyc} + \mathcal{L}_{seg}$)后,模型在下游分割任务中的表现大幅提升。与基线相比,SAR 到光学的 PA 指标升至 0.4480,表明显式的语义约束能有效指导生成器恢复地物的语义类别信息,减少语义模糊。为验证语义引导对比学习损失 \mathcal{L}_{sc} 的有效性,设置了将其替换为原始 CUT 方法中对比损失 \mathcal{L}_{nce} 的实验。实验结果表明,使用传统对比学习损失($\mathcal{L}_{gan} + \mathcal{L}_{cyc} + \mathcal{L}_{nce}$)虽然在一定程度上提升了图像纹理质量,但在下游任务的 PA 指标(0.3848)上明显低于使用语义对比损失($\mathcal{L}_{gan} + \mathcal{L}_{cyc} + \mathcal{L}_{sc}$)取得的 PA 指标(0.4923)。验证了语义引导策略通过类别一致性筛选正负样本,能更精准地实现特征对齐。

本文提出的完整模型(\mathcal{L}_{all})在 PSNR、SSIM 及 MAE 等图像质量指标上均取得了最优结果(PSNR 达到 21.1299),同时在下流分割任务中保持了最高的像素准确率(PA 为 0.4936)。尽管在光学到 SAR

的反向转换中,去除语义引导对比损失($\mathcal{L}_{gan} + \mathcal{L}_{cyc} + \mathcal{L}_{seg}$)的模型在PSNR上略有优势,但完整模型在双向转换任务中展现出了最佳的综合性能平衡,验证了融合循环生成结构与语义引导对比学习框架的有效性。

为了验证本文提出的联合优化框架的收敛性,图7展示了最优超参设置下($\lambda_1=10, \lambda_2=0.5, \lambda_3=1, \lambda_4=1, \lambda_5=1$)各项损失在训练过程中的变化趋势。其中, \mathcal{L}_{cyc} (红线)初始值最高,且随训练推进呈现持续且稳定的下降趋势,表明反向图像生成的一致性逐渐收敛。 \mathcal{L}_{seg} (橙线)与 \mathcal{L}_{sc} (绿线)均呈现显著下降趋势,并在约60 Epoch后收敛至稳定水平,说明模型成功学习并优化了跨模态语义一致性和特征对齐能力。 \mathcal{L}_{gan} (蓝线)在整个训练过程中保持在相对稳定的区间,没有出现消失或激增。这符合生成对抗网络博弈均衡的典型特征,说明生成器与判别器的性能协同提升并趋于稳定。以上所有损失项均能良好收敛,验证了本文所提联合优化框架的训练稳定性与整体有效性。

同时,为验证不同语义分割模型、分割损失权重、语义引导对比损失权重及图像块大小对本文方法的影响,在WHU-OPT-SAR数据集上进行了消融实验。表5展示了本文方法基于三类经典语义分割模型(Unet、DeepLabV3、Segformer)的定量实验结果,其中PA(real)表示各模型对真实SAR或光学图像的分割像素准确率。实验结果表明,在真实SAR图像上:三个模型的分割效果差异较小,其中DeepLabV3表现最优。本文方法基于DeepLabV3在图像生成质量和下游语义分割结果上均达到最优。在真实光学图像上:最优模型Segformer与最弱模型Unet的分割精度相差约6.68%。然而,基于Segformer的本文方法并未在图像生成质量和下游分割结果上取得最优;相对的,基于次优模型DeepLabV3的本文方法在PSNR和下游分割结果上均最优。综上,不同分割模型及其结果对本文方法有一定影响,但当分割效果达到一定阈值后,该影响趋于有限。因此,选择较易实现且常用的DeepLabV3模型为本文方法提供语义分割结果。

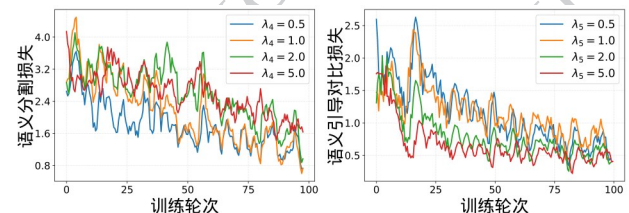
表6展示了不同语义分割损失权重(λ_4)下的定量实验结果。在光学到SAR转换任务中,本文方法在 $\lambda_4=1$ 时表现出最佳的综合性能,其PSNR、MAE和PA指标均达到最优,SSIM亦保持次优水平;权重过

低(0.5)或过高(5)均会导致各项指标出现不同程度的下降。在SAR到光学转换任务中,本文方法在 $\lambda_4=1$ 时生成图像的PSNR和MAE达到最优;虽然将权重增加至2能略微提升SSIM和PA,但代价是PSNR值出现了明显下降(从21.1299降至19.9154)。为在双向转换任务中取得最佳的性能均衡,最终将语义分割权重设置为1。

表7展示了不同语义引导对比损失权重(λ_5)下的定量实验结果。在光学到SAR转换任务中,模型在 $\lambda_5=1$ 时取得最佳性能,权重过低(0.5)或过高(5)均会对图像生成质量及下游语义分割任务的性能造成损害。在SAR到光学转换任务中, $\lambda_5=1$ 时PSNR值达到最优;继续增加 λ_5 会改善SSIM和PA,但代价是PSNR和MAE会更差。为在双向转换任务中取得最佳的性能均衡与稳定性,最终将语义引导对比损失权重设置为1。

为了进一步验证语义分割损失和语义引导对比损失在不同权重设置下的训练稳定性,图8展示了不同权重设置下两个训练损失收敛过程。实验结果表明,语义分割损失和语义引导对比损失在不同的权重设置下均能取得较好的训练收敛效果,表明这两个损失函数与提出方法有很好的适配性。

表8展示了语义引导对比学习模型在不同图像块大小配置下的定量实验结果。当图像块大小设为 1×1 时,模型能够利用语义信息最精准地筛选正负样本,从而在光学到SAR转换任务上取得了最优的图像生成质量与语义分割性能,并在SAR到光学转换任务上取得了最优的PSNR与MAE值及次优的PA值。该结果验证了 1×1 配置的有效性。结合传统对比学习也普遍采用此尺寸的实践,本研究最终将图像块大小确定为 1×1 。



图中曲线已进行移动平均平滑处理,窗口大小设置为5

图8 语义分割损失和语义引导对比损失在不同权重设置下训练的变化曲线

Fig. 8 Training curves of the semantic segmentation loss and the semantic-guided contrastive loss under different weight settings

表 4 消融实验结果

Table 4 Module ablation experiment results

损失组成	SAR→OPT				OPT→SAR			
	PSNR↑	SSIM↑	MAE↓	PA↑	PSNR↑	SSIM↑	MAE↓	PA↑
\mathcal{L}_{gan}	17.4187	0.3579	0.1384	0.1165	11.4415	0.0214	0.2053	0.0645
$\mathcal{L}_{\text{gan}}+\mathcal{L}_{\text{eye}}$	18.2425	0.3629	0.1185	0.2701	11.8556	0.0290	0.1931	0.3026
$\mathcal{L}_{\text{gan}}+\mathcal{L}_{\text{eye}}+\mathcal{L}_{\text{seg}}$	20.8765	<u>0.3720</u>	<u>0.0806</u>	0.4480	13.2284	0.0578	0.1658	<u>0.4694</u>
$\mathcal{L}_{\text{gan}}+\mathcal{L}_{\text{eye}}+\mathcal{L}_{\text{sc}}$	20.5921	0.3704	0.0831	<u>0.4923</u>	12.0768	0.0401	0.1873	0.4373
$\mathcal{L}_{\text{gan}}+\mathcal{L}_{\text{eye}}+\mathcal{L}_{\text{ncc}}$	19.2422	0.3679	0.1026	0.3848	12.0383	0.0440	0.1867	0.4128
\mathcal{L}_{all}	21.1299	0.3793	0.0765	0.4936	<u>12.9165</u>	<u>0.0562</u>	<u>0.1705</u>	0.5061

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

表 5 不同分割模型的消融实验结果

Table 5 Ablation experiment results of different segmentation models

分割模型	SAR→OPT					OPT→SAR				
	PSNR↑	SSIM↑	MAE↓	PA↑	PA(real)↑	PSNR↑	SSIM↑	MAE↓	PA↑	PA(real)↑
Unet(Ronneberger 等, 2015)	<u>20.6731</u>	0.3523	0.0744	0.3474	0.6861	<u>12.3943</u>	<u>0.0532</u>	<u>0.1811</u>	<u>0.3774</u>	0.7822
Segformer(Xie 等, 2021)	19.5814	0.4463	0.0955	<u>0.3991</u>	0.7319	12.0282	0.0325	0.1892	0.3090	<u>0.7853</u>
DeepLabV3(Chen 等, 2017)	21.1299	<u>0.3793</u>	<u>0.0765</u>	0.4936	<u>0.7096</u>	12.9165	0.0562	0.1705	0.5061	0.7868

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

表 6 语义分割权重 λ_4 的消融实验结果Table 6 Ablation results of the semantic segmentation weight λ_4

λ_4	SAR→OPT				OPT→SAR			
	PSNR↑	SSIM↑	MAE↓	PA↑	PSNR↑	SSIM↑	MAE↓	PA↑
0.5	<u>20.4261</u>	0.3891	0.0865	0.4861	12.1582	<u>0.0463</u>	0.1857	0.3373
1	21.1299	0.3793	0.0765	<u>0.4936</u>	12.9165	<u>0.0562</u>	0.1705	0.5061
2	19.9154	0.4032	0.0891	0.5110	<u>12.7364</u>	0.0565	0.1732	0.3782
5	19.7951	<u>0.4009</u>	<u>0.0820</u>	0.4767	12.6271	0.0544	<u>0.1720</u>	<u>0.3783</u>

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

3 结论

本文提出了一种语义分割引导对比学习的 SAR

到光学图像转换方法,设计了融合循环生成与对比学习的联合优化框架,有效解决了传统对比学习在遥感图像转换中的特征同质化问题。该方法设计了基于语义分割的类别一致性正负样本筛选策略,构

表7 语义引导对比损失权重 λ_s 的消融实验结果Table 7 Ablation results of the semantic-guided contrastive loss weight λ_s

λ_s	SAR→OPT				OPT→SAR			
	PSNR↑	SSIM↑	MAE↓	PA↑	PSNR↑	SSIM↑	MAE↓	PA↑
0.5	<u>21.0849</u>	0.4098	0.0740	0.4332	<u>12.2567</u>	<u>0.0508</u>	<u>0.1845</u>	<u>0.3837</u>
1	21.1299	0.3793	0.0765	0.4936	12.9165	0.0562	0.1705	0.5061
2	20.2101	0.4239	0.0847	0.5206	12.1768	0.0399	0.1862	0.3095
5	21.0287	<u>0.4176</u>	<u>0.0747</u>	<u>0.5124</u>	11.9662	0.0387	0.1901	0.3232

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

表8 图像块不同大小的消融实验结果

Table 8 Ablation experiment results with different patch sizes

块大小	SAR→OPT				OPT→SAR			
	PSNR↑	SSIM↑	MAE↓	PA↑	PSNR↑	SSIM↑	MAE↓	PA↑
5×5	<u>20.6379</u>	<u>0.4331</u>	<u>0.0785</u>	0.5253	<u>11.9701</u>	<u>0.0343</u>	<u>0.1909</u>	<u>0.3267</u>
3×3	19.9700	0.4350	0.0880	0.4505	11.9540	0.0340	0.1911	0.2945
1×1	21.1299	0.3793	0.0765	<u>0.4936</u>	12.9165	0.0562	0.1705	0.5061

注:加粗字体为每列最优值,下划线字体为次优;“↑”表示数值越大越好,“↓”表示数值越小越好。

建了语义引导的对比学习模块,有效提升了跨模态特征对齐的准确性;同时,引入循环语义分割损失,约束了生成图像的结构、纹理与语义的一致性,从而提升了生成图像的质量与实用性。在 WHU-OPT-SAR 和 DDHRNet 两个公开数据集上的实验结果表明,与相关方法相比,本文方法不仅能生成高质量、高保真度的跨模态图像,还能显著提升生成图像在语义分割和图像匹配等下游任务中的性能。然而,本文方法依赖于预训练的语义分割模型来提供先验知识,这限制了方法在缺乏语义标注数据场景下的应用,且分割模型的精度会直接影响最终的转换效果。同时,当前方法中正负样本的构建严格依赖于预定义的类别标签,这种强监督设定在跨场景或完全无标注的环境中适应性有限,难以支持真正的无监督泛化需求。因此,下一步将研究如何在无需预训练语义分割模型的情况下,实现自监督的语义引导对比学习,并探索不依赖显式类别标签的正负样本构建机制,以增强方法的通用性和鲁棒性。

参考文献 (References)

- Bai X, Pu X, Xu F. 2023. Conditional diffusion for SAR to optical image translation. *IEEE Geoscience and Remote Sensing Letters*, 21: 1-5 [DOI: 10.1109/LGRS.2023.3337143]
- Badrinarayanan V, Kendall A, Cipolla R. 2017. SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39 (12): 2481-2495 [DOI: 10.1109/TPAMI.2016.2644615]
- Chen L C, Papandreou G, Schroff F, Adam H. 2017. Rethinking atrous convolution for semantic image segmentation [EB/OL]. [2025-09-03].
<https://arxiv.org/pdf/1706.05587.pdf>
- Guo Z, Zhang Z, Cai Q, Liu J, Fan Y, Mei S. 2024. MS-GAN: Learn to memorize scene for unpaired SAR-to-optical image translation. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 17: 11467-11484 [DOI: 10.1109/JSTARS.2024.3411691]
- Han J L, Shoeiby M, Petersson L, Armin M A. 2021. Dual contrastive learning for unsupervised image-to-image translation // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. Online: IEEE: 746-755 [DOI: 10.1109/CVPRW53098.2021.00084]
- Hu J M, Li Y L, Zhi X Y, Shi T J, Zhang W. 2025. Complementarity-

- aware feature fusion for aircraft detection via unpaired Opt2SAR image translation. *IEEE Transactions on Geoscience and Remote Sensing*, 63: 1-19 [DOI:10.1109/TGRS.2025.3578876]
- Hu X Q, Zhou X Y, Huang Q S, Shi Z Y, Sun L, Li Q L. 2022. QS-Attn: query-selected attention for contrastive learning in I2I translation // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. New Orleans, USA: IEEE: 18270 - 18279 [DOI:10.1109/CVPR52688.2022.01775]
- Jung C Y, Kwon G Y, Ye J C. 2022. Exploring patch-wise semantic relation for contrastive learning in image-to-image translation tasks // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. New Orleans, USA: IEEE: 18239 - 18248 [DOI:10.1109/CVPR52688.2022.01772]
- Jiang W D, Sun Y L, Lei L, Kuang G Y, Ji K F. 2025. AdaptVFM-RSCD: Advancing remote sensing change detection from binary to semantic with SAM and CLIP. *ISPRS Journal of Photogrammetry and Remote Sensing*, 230: 304-317 [DOI: 10.1016/j.isprsjprs.2025.09.010]
- Kang M, Park J. 2020. ContraGAN: contrastive learning for conditional image generation // *Proceedings of the 34th International Conference on Neural Information Processing Systems*. Vancouver, Canada: NIPS: 21357 - 21369
- Li X, Zhang G, Cui H, Wang S, Li X, Chen Y, Li Z, Zhang L. 2022. MCANet: A joint semantic segmentation framework of optical and SAR images for land use classification. *International Journal of Applied Earth Observation and Geoinformation*, 106: 1-13 [DOI: 10.1016/j.jag.2021.102638]
- Park T, Efros A A, Zhang R, Zhu J Y. 2020. Contrastive learning for unpaired image-to-image translation // *Proceedings of the European Conference on Computer Vision (ECCV)*, Part IX. Online: Springer: 319 - 345 [DOI:10.1007/978-3-030-58545-7_19]
- Ren B, Ma S, Hou B, Hong D, Chanussot J, Wang J, Jiao L. 2022. A dual-stream high resolution network: Deep fusion of GF-2 and GF-3 data for land cover classification. *International Journal of Applied Earth Observation and Geoinformation*, 112: 1-15 [DOI: 10.1016/j.jag.2022.102896]
- Ronneberger O, Fischer P, Brox T. 2015. U-Net: Convolutional networks for biomedical image segmentation // *Proceedings of the 18th International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Part III. Munich, Germany: Springer: 234 - 241 [DOI: 10.1007/978-3-319-24574-4_28]
- Torbunov D, Huang Y, Yu H W, Huang J, Yoo S, Lin M F, Viren B, Ren Y H. 2023. UVCGAN: UNet Vision Transformer cycle-consistent GAN for unpaired image-to-image translation // *Proceedings of the 23rd IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. Waikoloa, USA: IEEE: 702 - 712 [DOI:10.1109/WACV56688.2023.00077]
- Wu L W, Sun R, Kan J S and Gao J. 2020. Double dual generative adversarial networks for cross-age sketch-to-photo translation. *Journal of Image and Graphics*, 25(4): 732-744 (吴柳玮, 孙锐, 阚俊松, 高隽. 2020. 双重对偶生成对抗网络的跨年龄素描—照片转换. *中国图象图形学报*, 25(4): 732-744) [DOI: 10.11834/jig.190329]
- Wang W L, Zhou W G, Bao J M, Chen D, Li H Q. 2021. Instance-wise hard negative example generation for contrastive learning in unpaired image-to-image translation // *Proceedings of the 18th IEEE/CVF International Conference on Computer Vision (ICCV)*. Montreal, Canada: IEEE: 14000 - 14009 [DOI: 10.1109/ICCV48922.2021.01376]
- Wang P, Chen Y K, Huang B, Zhu D Y, Lu T W, Dalla Mura M, Chanussot J. 2025. MT_GAN: A SAR-to-optical image translation method for cloud removal. *ISPRS Journal of Photogrammetry and Remote Sensing*, 225: 180 - 195 [DOI:10.1016/j.isprsjprs.2025.04.011]
- Xie E Z, Wang W H, Yu Z D, Anandkumar A, Alvarez J M, Luo P. 2021. SegFormer: Simple and efficient design for semantic segmentation with Transformers // *Proceedings of the 35th Annual Conference on Neural Information Processing Systems (NeurIPS)*. Online: NIPS: 12077-12090
- Yu P L, Shi Q and Wang H. 2021. Infrared-to-visible image translation based on parallel generator network. *Journal of Image and Graphics*, 26(10): 2346-2356 (余佩伦, 施佳, 王晗. 2021. 并行生成网络的红外—可见光图像转换. *中国图象图形学报*, 26(10): 2346-2356) [DOI:10.11834/jig.200113]
- Yang X, Wang Z, Zhao J, Yang D. 2022. FG-GAN: A fine-grained generative adversarial network for unsupervised SAR-to-optical image translation. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1-11 [DOI:10.1109/TGRS.2022.3165371]
- Yagmur I C, Ates H F, Gunturk B K. 2024. XPoint: A Self-Supervised Visual-State-Space based Architecture for Multispectral Image Registration[EB/OL].[2025-09-03]. <https://arxiv.org/pdf/2411.07430.pdf>
- Zhan F N, Zhang J H, Yu Y C, Wu R L, Lu S J. 2022. Modulated contrast for versatile image synthesis // *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. New Orleans, USA: IEEE: 18259 - 18269 [DOI: 10.1109/CVPR52688.2022.01774]
- Zhu J Y, Park T, Isola P, Efros A A. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks // *Proceedings of the 16th IEEE International Conference on Computer Vision (ICCV)*. Venice, Italy: IEEE: 2242 - 2251 [DOI:10.1109/ICCV.2017.244]

作者简介

杜文亮,男,讲师,主要研究方向为遥感图像处理。E-mail: wldu@cumt.edu.cn

郭波,男,硕士研究生,主要研究方向为遥感图像转换。E-mail: boguo@cumt.edu.cn

赵佳琦,男,副教授,主要研究方向为人工智能、计算机视觉。E-mail: jiaqizhao@cumt.edu.cn

姚睿,男,教授,主要研究方向为人工智能、计算机视觉。E-mail: ruiyao@cumt.edu.cn

周勇,通信作者,男,教授,主要研究方向为机器学习、人工智能。E-mail: yzhou@cumt.edu.cn